

A Semantic Continuity Based Analysis of Topic Lifecycle on Social Networks

Kuntal Dey^{1(✉)}, Saroj Kaushik², Kritika Garg³, and Ritvik Shrivastava⁴

¹ IBM Research, New Delhi, India

kuntadey@in.ibm.com

² Indian Institute of Technology, Delhi, India

saroj@cse.iitd.ac.in

³ Ch. Brahm Prakash Government Engineering College, New Delhi, India

kgarg.kritika@gmail.com

⁴ Netaji Subhas Institute of Technology, New Delhi, India

ritviks.it@nsit.net.in

Abstract. Analyzing the lifecycle of topics, that are present in user-generated text content, has emerged as a mainstream topic of social network research. The literature presently identifies topics on Twitter, a prominent online social network, as either individual hashtags, or a burst of keywords within a short span of time, or as latent concept spaces obtained from sophisticated text analysis mechanisms, such as Latent Dirichlet Allocation (LDA). The first and second approaches fail to recognize that topics do not restrict themselves to individual hashtags and are likely to span across (semantically related) keywords, while the third does not capture the user’s intended topics expressed via hashtags. In the current paper, we propose a novel methodology that addresses these shortcomings. We jointly utilize the temporal concurrency of the hashtags contained in given tweets and the latent concept space addressed by the tweet content, to identify groups of hashtags representing concept space—a “topic”—addressed by many tweets. A given topic, thus, is represented by a different set of representative hashtags at different times; the usage rate of the different hashtags change such that some hashtags gain prominence over others over time. Unlike the literature, where lifecycle analysis of one topic typically comprises of analyzing one hashtag, we analyze and characterize the lifecycle of a topic as a combination of multiple semantically and temporally related hashtags. We derive novel insights about lifecycle of topics: the inception and continuity of the topics over time (expressed over different hashtags), and how topics morph over hashtags, from one set of hashtags to another, before eventually dying down.

1 Introduction

Twitter, a hotbed of user-generated content, has been a platform of intense research focus for social network researchers from the angle of information content analysis. One key research area has been centered around the topics of

user-generated content, attempting to identify the topics of user interest, understanding the lifecycle of these topics—how these topics are born, spread over and lead to user attention over time, and eventually die down, and analyzing the information that diffuse with these topics over the social network users. This is important for obtaining a better understanding of information diffusion dynamics on online social networks.

1.1 Related Work (Background) and Motivation

The literature for topic identification on Twitter has followed three different approaches. In the first approach, the hashtags that are part of the tweet messages, have been treated as the topic as-is. Works, such as [6], have used this approach to associate topics with tweets. In the second approach, a burst of keywords in a short span of time are identified, and each bursting keyword is treated as a topic. Works, such as [4, 5, 11], have used this approach. And in the third approach, the latent semantic concepts of given tweets—often identified with sophisticated text-to-topic assignment techniques such as Latent Diriclet Allocation (LDA) [3]—are treated as topics, and the tweets that address these spaces are said to belong to these topics. Works, such as [10], have followed this approach. A detailed survey of the works, that follow these different bodies of work, can be found in [7].

However, these approaches suffer from inherent limitations. The first two approaches miss out on the latent semantic concept space addressed by the content, since they simply examine the keywords and hashtags instead of the overall content space. Thus, for example, the first approach would not be able to identify that tweets containing hashtags *#mj*, *#michaeljackson*, *#jackson* and *#m.jackson* potentially address the same topic, or at least, need an investigation on whether these intend to address the same topic or not. Similarly, the second approach would not attempt to identify whether the intent of the keywords *mj*, *michaeljackson*, *jackson* and *m.jackson* are potentially the same or not. The third approach attempts to identify the semantic concept space, but without accounting for the explicit user intent that a hashtag is intended for, as the studies conducted in the literature typically consider the content without giving any emphasis to the user intent (expressed in form of hashtags).

Twitter hashtag disambiguation has emerged as a topic of research interest. In one of the early works, Yang and Leskovec [17] detected similar shapes (distributions) of usage of given Twitter hashtags, using K-Spectral Centroid (KSC) clustering. However, this work neither inspected the aspect of temporal overlap of different hashtags (contemporary tweets using the same hashtag), nor did it attempt to explore the semantic concept space covered by the tweets under consideration. In a subsequent work, [16] propose SAX, a temporal sense clustering algorithm based on the idea that semantically related hashtags have similar and synchronous usage patterns. The SAX framework, by its nature, overcomes one primary shortcoming of KSC, by incorporating the temporal overlap of hashtags. However, it does not take the semantics of the content into account, which leads to clustering of topically unrelated tweets also.

None of the hashtag disambiguation methodologies account for topics morphing (evolving), from being primarily constituted of one group of hashtags to other groups of hashtags, over time. In addition, none of the topic identification and topic lifecycle analysis works, take hashtag disambiguation and topic morphing jointly into consideration. In a work that performs a thorough analysis of Twitter topic lifecycles, Ardon *et al.* [2] used hashtag as an indicator of topic (in their jargon, “event”, which is in philosophy similar to what we call “topic” here in absence of any unified jargon or terminology in the research community). Other works, such as [8, 13, 15], also use hashtag and LDA based methods for identifying topics, and analyzing their spatio-temporal evolution. However, none of these works attempt to investigate the lifecycle of topics, from the perspective of content semantics, along with the temporal overlaps. In other works, while in principle it is feasible to extend the semantic concept space based approach for hashtag disambiguation and thus better understand topic lifecycles and information diffusion, no such study has been carried out in the literature on this front.

1.2 Contributions of Our Work

We hypothesize that, hashtags, and topics derived thereof, have the following characteristics.

- Multiple hashtags are used by different Twitter users towards the same topic. Further, since the topics are the same, there is some temporal overlap amongst the time period that these hashtags are used. If such semantically and temporally related hashtags are better identified, the characteristics of information networks can be better identified.
- While the usage of some hashtags reduce over time, that does not imply the end of the topic. Other related hashtags, that address the same topic, could become more popular and “replace” the hashtags that were prominent for the topics earlier. Thus, the discussion lifecycle of the underlying topic could be longer than just the lifespan of one or a few hashtags. As of today, the literature tends to treat a hashtag or a bursty keyword as a topic; however, understanding how topics morph over time, moving from one hashtag to another while addressing similar concepts, and eventually die down rather than morphing, would add a novel understanding towards the temporal continuity of discussion topics on Twitter.

We create a timeline for the hashtags, that tracks the usage frequency of the hashtags at each given time. We identify the latent semantic concept space addressed by each tweet using LDA, and assign each tweet to the main latent concept it addresses. This, in turn, leads to a mapping between the latent concepts and hashtags that are part of given tweets belonging to that latent concept, forming per-concept hashtag clusters. We treat each (semantically and temporally related) hashtag cluster obtained as a topic, and treat all the corresponding tweets to belong to this topic. We inspect the temporal manifestation of the

hashtags. We define the (temporal) lifecycle of a topic as the combined (temporal) lifecycle of the individual hashtags, and the intensity of a topic at a given time period (on the timeline) as the total (concurrent) occurrence of the individual hashtags within that time period. We observe the dominance of certain hashtags over others, at certain times, within given topics, as well as, evolution of the dominance levels of the hashtags. Using our methodology of deriving semantically related topics that witness continuation of discussions morphing over different hashtags, we provide a study of some key characteristic attributes of the corresponding social network graph. We demonstrate the effectiveness of our approach using around 20–30% of eighteen days of Twitter data.

2 Our Approach

In this section, we present the details of our approach. As mentioned earlier in Sect. 1, the primary objectives of our work include (a) identifying topics as a set of semantically related hashtags that associate with tweets conveying semantically similar (overlapping) concepts (we call these as LDA-concepts) and are temporally related, and (b) subsequently analyzing the lifecycle of topics that are formed as a set of such semantically similar hashtags, which is a first-of-its-kind analysis in the literature: the earlier works, such as [2], focused on analyzing the lifecycle of individual hashtags. We also explore the properties of the social connection graph of users built around each such topic, to characterize the topic-specific user graph thus formed.

2.1 Identification of Related Hashtags

The input to our system is a collection of tweets, collected from Twitter, that comprise of at least one hashtag. Our aim in this step is to create clusters of hashtags, such that each cluster consists of semantically related tweets.

2.1.1 Topic Modeling Using LDA

We perform LDA-based topic modeling to identify the latent semantic concept space contained in the tweets. to avoid confusion of terminology (specifically the overloaded term *topic*), going forward, we shall refer to the LDA-obtained topics as *LDA-concepts* while the final hashtag clusters we obtain with semantic and temporal relationships as *topics*. We perform this over the following pair of steps.

First, we create a document as a set union (concatenation) of all the tweets in the input data, ignoring the hashtags. That is, for a given set of tweets $T = \{t_1, t_2, \dots, t_n\}$, that consist of hashtags $H = \{h_1, h_2, \dots, h_m\}$, we obtain a document D as

$$D = \bigcup_{i=1}^n t_i - \bigcup_{j=1}^m h_j \quad (1)$$

Next, we perform LDA-concept modeling using the document D as input, and learn a set of LDA-concepts $Z = \{z_1, z_2, \dots, z_l\}$. Note that, LDA [3] is traditionally modeled as a joint distribution as follows:

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D}) = \prod_{i=1}^K p(\beta_i) \cdot \prod_{d=1}^D p(\theta_d) \cdot \prod_{n=1}^N p(z_{d,n} | \theta_d) p(w_{d,n} | \beta_{1:K}, z_{d,n}) \quad (2)$$

Here, θ_d are the LDA-concept proportions for document d , $\theta_{d,k}$ is the LDA-concept proportion for LDA-concept k in document d , z_d are the LDA-concept assignments for document d , $\beta_{1:K}$ are the LDA-concepts where each β_k is a distribution over the given vocabulary, w_d are the observed words for document d , and $z_{d,n}$ is the LDA-concept assignment for word n in document d . The two steps above complete the latent semantic concept space learning, and extracting the LDA-concepts globally present in the documents.

2.1.2 Semantically and Temporally Related Hashtag Cluster Creation

Semantically related hashtag cluster creation

Following a non-overlapping clustering scheme, we associate each tweet with a single LDA-concept. In order to perform this, we compute the probability of each given tweet to belong to each of the LDA-concepts $\{z_1, z_2, \dots, z_l\}$ obtained in Eq. 2. We rank these probabilities, and assign the tweet to the LDA-concept z_k for which the probability is the highest. Each hashtag h_i belonging to a given tweet t_j is thus assigned to the LDA-concept z_k .

Note that, since the same hashtag can be used in different tweets under different semantic concept contexts, a hashtag need not necessarily uniquely belong to a single LDA-concept. For instance, the hashtag *#president* can be used in the context of a country, or a club, or an organization, and the corresponding tweets will probably belong to different LDA-concepts. Thus, the hashtag will belong to different clusters, under different context of usage. Our model, by its inherent design, accounts for such contextual usages.

Temporally relating hashtags

Hashtags that need to be clustered together for identifying the temporal continuation of a given semantic topic, need to be temporally related also, along with their semantic relationship under the context of given content that has been detected above. Allen [1] created an exhaustive list of temporal relationships that can exist between a pair of time periods. This includes *overlap*: part of event A and event B co-occur, *meets*: event A starts as soon as event B stops, and *disjoint*: event A and event B share no common time point. In our case, an event is an instance of a tweet using a given hashtag. We create a time series of the individual hashtags, as well as the hashtag clusters. At each slot in the time series (set to a given time duration, e.g., 1 h etc.), the number of times each hashtag is used, is computed. Two hashtags h_i and h_j will be temporally related

as long as they either satisfy the *overlaps* relationship, or if there exists one or more hashtags h_k , such that, h_i is temporally related to h_k , and h_k *overlaps* h_j . Two hashtags h_i and h_j are temporally unrelated if $\nexists h_k$ such that h_i is temporally related to h_k , and, h_k *overlaps* h_j . The *temporally related* relationship is recursive in nature, and can be expressed as

$$h_i \odot h_j \implies \left((\exists h_k) h_i \odot h_k \right) \cap (h_k \odot h_j) \quad (3)$$

where \odot denotes the *temporally related* relationship and \odot denotes the *overlaps* relationship. We finalize our topics, defined as hashtag clusters, such that each hashtag cluster consists of hashtags that are both semantically and temporally related.

As an example, at the end of the process, hashtags such as $\{\#tennis, \#federer, \#rogerfederer, \#roger\}$ *etc.* are expected to be together in one cluster together, while hashtags such as $\{\#politics, \#trump, \#donlandtrump, \#donald\}$ *etc.* are expected to be together another cluster together.

2.2 Topic Lifecycle Analysis

Ardon *et al.* [2] observe that events have five phases in their lifecycles: pre-growth phase, growth phase, peak phase, decay phase and post-phase. In pre-growth phase, an event is introduced into a network. Early adopters talk about it in the growth phase. In the peak and decay phases, and “early majority” and a “late majority” talk about it. Finally, laggards talk about it in the post-phase. They treat hashtags as topics for tweets containing hashtags, and use a NLP tool to identify places, entities *etc.* and assign these as tags (in turn, these tags become topics).

On the other hand, we do not attempt to solve the hashtag assignment problem to the posts that do not contain hashtags (as it is an independent problem altogether). Using only the tweets that contain hashtags, we attempt to perform analysis of lifecycle of the topics. Note that, the hashtag-cluster based topic identification ensures that each tweet under consideration will comprise of at least one hashtag. Since the topics formed by our approach comprise of a composition of hashtags, we investigate two main aspects.

Dominant hashtag detection and topic morphing

Let for a given topic z_k , $H = \{h_{1,k}, h_2, \dots, h_m\}$ be the set of hashtags used at each given time slot (within the time series). If c is a function that counts the number of times each hashtag h_i was used, then, the dominant hashtag for the given timeslot is defined as $h_x = \forall(i)(\max(c(h_i)))$. Essentially, it is the most frequently used hashtag of a given cluster within a given time slot. Note that, for a given topic, each individual hashtag, including each hashtag that has been dominant in one or more time slots, tends to follow a lifecycle similar to what was observed for events by Ardon *et al.* [2]. However, as soon as in a new (next) slot a different hashtag replaces an earlier dominant hashtag and becomes the dominant hashtag of its cluster in the slot, the intensity of the topic covered

by the cluster remains high. Thus, while each individual hashtag has a smaller lifecycle, the dominant hashtag ensures that the “most representative hashtag” of each specific slot shows the presence of the topic at a sufficient level. This also shows how a topic morphs, via one set of hashtags (or, from one single hashtag) to another.

Topic intensity detection

Let for a given topic z_k , $H = \{h_{1,k}, h_{2,k}, \dots, h_{m,k}\}$ be the set of hashtags used at each given time slot (within the time series). If c is a function that counts the number of times each hashtag h_i was used, then, the topic intensity at the given timeslot is defined as $k_{h_x} = \sum_{i=1}^m (c(h_i))$. Essentially, topic intensity as computed as the total usage of all the hashtags of a given cluster within a given time slot, giving the cumulative (total) presence of a topic within the time slot.

Note that, in our experiments, we perform a complete study of the lifecycle of the topics (hashtag clusters), that include observing the different parts of the bigger topic (the combination of the hashtags) as well as the individual hashtags, and contrast between some of the observations for individual hashtags versus the hashtag clusters. We also make notes for characterizing the early, mid and late parts of the lifecycles of individual hashtags, and observing how topics morph from one set of dominant hashtags at one time to a different set, as well as, investigate whether and how topics represented by a multiplicity of hashtags eventually die.

2.3 Graph Characteristics

We finally attempt to explore the social network subgraphs obtained for each of the topics, by examining metrics such as degree distribution, graph diameter, and strongly and weakly connected components. We omit the definitions of these metrics, as these are fundamental to graph theory, and are well-known enough to expect reader’s familiarity to their definitions.

3 Experiments

Dataset Description

Our experiments are conducted using a popular Twitter dataset¹ made available by Yang and Leskovec [17]. The dataset claims to have 20–30% of the entire set of tweets within the given period of time. We use the entire set of tweets of the last 18 days of June 2009, and retain the social connections of the users using the corresponding social networks made available² by Kwak *et al.* [9]. Akin to the approach of Ardon *et al.* [2], we perform experiments by retaining the most highly used 1,000 hashtags within the given period of time, and the tweets thereof. We retain the users that posted these tweets, and the social connections

¹ <https://snap.stanford.edu/data/twitter7.html>.

² <http://an.kaist.ac.kr/traces/WWW2010.html>.

Table 1. Description of Available Data. All the tweets are from June 2009.

Total num. of tweets	Num. tweets retained	Num. users retained	Num. edges retained	Avg. num. tweets/user	Avg. num. connections
18,572,084	1,339,272	117,701	4,973,218	11.38	42.25

among these users, to form the relevant social network subgraph. The statistics of the dataset are presented in Table 1.

Experimental Setup

We conduct our experiments on the given data, following the steps delineated in Sect. 2. We find the LDA-concepts, and assign clusters to the input tweets, using MALLET [12]. We repeat our experiments at different granularities of topics, namely by detecting 50, 100, 150, 200, 250, 300, 350, 400, 450 and 500 topics from the given 1000 hashtags, to observe the quality of results at different possible granularities. We subsequently create the timelines for each individual hashtag, as well as each topic comprising of a cluster of hashtags, and conduct the rest of the experiments on this.

Topic Lifecycles: Topic Intensity and Hashtag Dominance Characteristics

The intensity of a topic (represented as a hashtag cluster), and dominance of hashtags, are studied. Figure 1 shows the characteristics of topic intensities over time. It can be seen that, topics reach peaks at periodic intervals, and then reduce in terms of intensities as hashtags morph into others, and subsequently pick up.

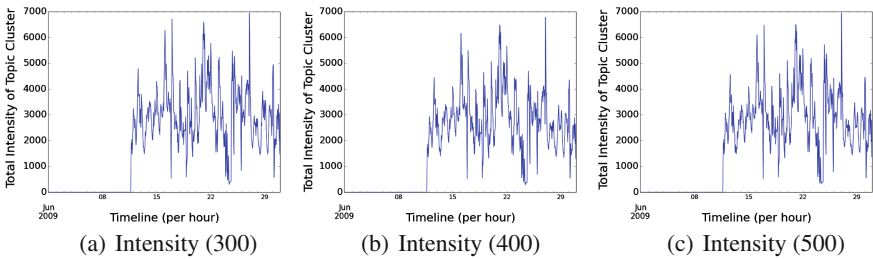


Fig. 1. Topic intensity over time. (a) Intensity (300) (b) Intensity (400) (c) Intensity (500)

Figure 2 depicts the varying patterns of dominance of different hashtags over time. We randomly pick a few representative samples for the purpose of illustration; however, we manually inspect and observe the core characteristics to be present in the rest of the results too. Clearly, while the rate of usage one hashtag tapers off over time, other hashtags take over and become dominant.

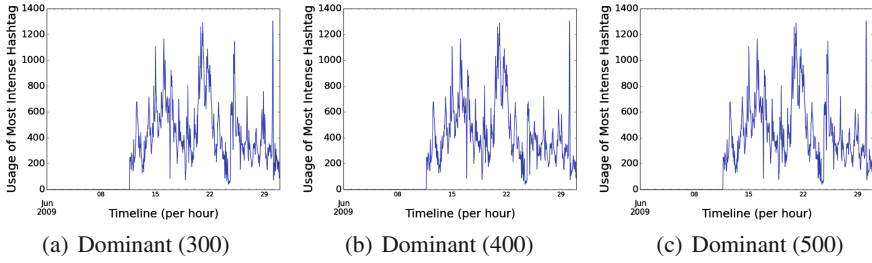


Fig. 2. Evolution of topics, where the dominant hashtags temporally evolve and replace one another. (a) Dominant (300) (b) Dominant (400) (c) Dominant (500)

Further, both Figs. 1 and 2 indicate a strong tendency of topics to morph over time, rather than die down in isolation, which is a novel insight obtained by this work.

Network Graph Characteristics

We study the characteristics of the social network graphs that are formed by retaining the edges where both the participating users have at least one tweet belonging to a given cluster (topic). Specifically, we find the graph diameter, and strongly and weakly connected components. We use the NETWORKX package available as part of the PYTHON, for computing the graph properties. The characteristics of the lifecycles of the topics (hashtag clusters) can be seen on Figs. 4, 5 and 6 for cluster sizes 300, 400 and 500 respectively. To characterize the input data, we explore its degree distribution (Fig. 3), and observe it to follow the long tail pattern that is well-established in the social network analysis literature by studies such as Nanavati *et al.* [14].

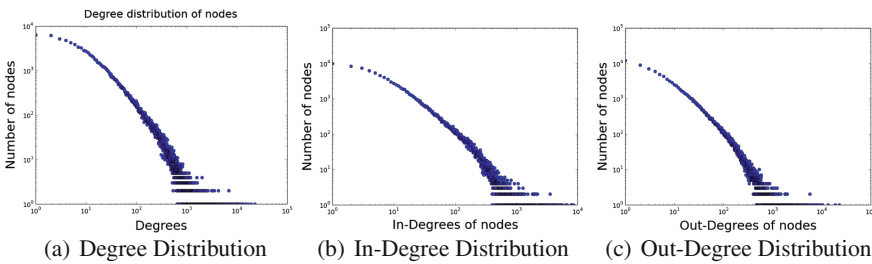


Fig. 3. Overall, in-degree and out-degree distribution. (a) Degree Distribution (b) In-Degree Distribution (c) Out-Degree Distribution

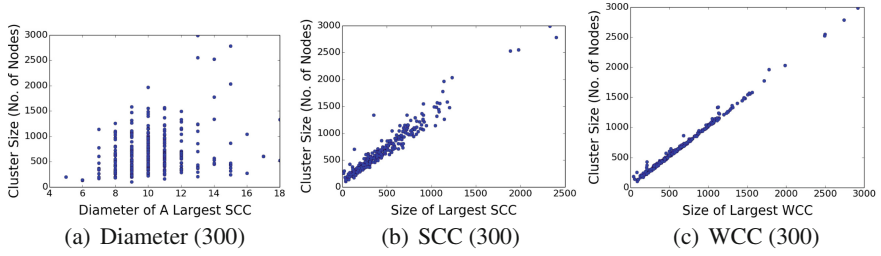


Fig. 4. Network graph characteristic properties for the 300-topic granularity. (a) Diameter (300) (b) SCC (300) (c) WCC (300)

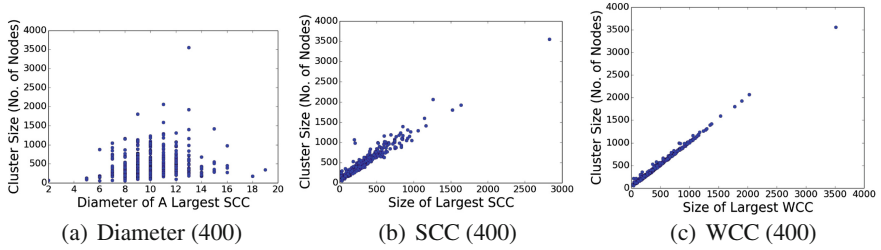


Fig. 5. Network graph characteristic properties for the 400-topic granularity. (a) Diameter (400) (b) SCC (400) (c) WCC (400)

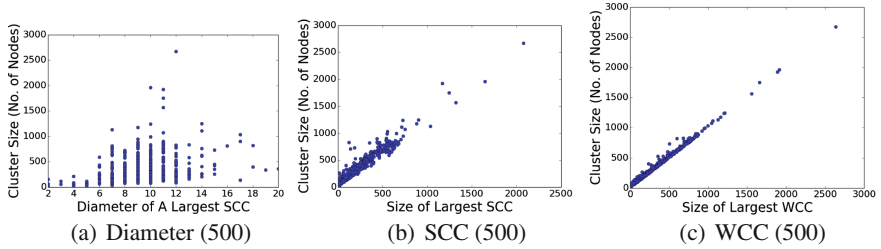


Fig. 6. Network graph characteristic properties for the 500-topic granularity. (a) Diameter (500) (b) SCC (500) (c) WCC (500)

4 Conclusion

In the current paper, we provided a novel analysis of topic lifecycles. In the process, we first proposed a framework that would use a higher-dimension latent semantic concept space for bringing tweets with similar concepts but different hashtags together. This enabled identifying clusters of related hashtags, that are temporally proximal or overlapping, and semantically related. Effectively, these hashtags are representatives of the semantic topic core to the group of the corresponding tweets. Subsequently, we analyzed the timeline of these hashtags, and characterized the lifecycle of the topics using metrics that have been user earlier

in the literature. We observed significant differences between the characteristics of topic lifecycle that has been given in the literature, and the characteristics that are obtained by our framework, in terms of (a) longevity, (b) user participation and (c) continuity of topics by morphing into different hashtags versus dying down completely. These differences can be attributed to our overall approach to the problem, which is different from the literature. We empirically demonstrated these differences, by comparing multiple graph characteristics we obtained using our hashtag clustering based topic lifecycle analysis framework, with that in the literature. Our work necessitates a revisit of information diffusion models on social networks, with a potential of updating our understanding of the dynamics of the diffusion with respect to topics and their lifecycles; we aim to do this in the future. In the future, we also aim to characterize the graph attributes of the individual clusters, as well as, cumulative graph attributes that span across clusters but can be socially interwoven.

References

1. Allen, J.F.: Maintaining knowledge about temporal intervals. *Commun. ACM* **26**(11), 832–843 (1983)
2. Ardon, S., Bagchi, A., Mahanti, A., Ruhela, A., Seth, A., Tripathy, R.M., Triukose, S.: Spatio-temporal and events based analysis of topic popularity in Twitter. In: *CIKM*, pp. 219–228. ACM (2013)
3. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. *J. Mach. Learn. Res.* **3**, 993–1022 (2003)
4. Cataldi, M., Di Caro, L., Schifanella, C.: Emerging topic detection on twitter based on temporal and social terms evaluation. In: *Tenth International Workshop on Multimedia Data Mining*, pp. 4. ACM (2010)
5. Cataldi, M., Schifanella, C., Candan, K.S., Sapino, M.L., Di Caro, L.: Cosena: a context-based search and navigation system. In: *International Conference on Management of Emergent Digital EcoSystems*, p. 33. ACM (2009)
6. Cunha, E., Magno, G., Comarela, G., Almeida, V., Gonçalves, M.A., Benevenuto, F.: Analyzing the dynamic evolution of hashtags on twitter: a language-based approach. In: *Proceedings of the Workshop on Languages in Social Media*, pp. 58–65. ACL (2011)
7. Dey, K., Kaushik, S., Subramaniam, L.V.: Literature survey on interplay of topics, information diffusion and connections on social networks (2017). arXiv preprint [arXiv:1706.00921](https://arxiv.org/abs/1706.00921)
8. Ifrim, G., Shi, B., Brigadir, I.: Event detection in twitter using aggressive filtering and hierarchical tweet clustering. In: *SNOW-DC@ WWW*, pp. 33–40 (2014)
9. Kwak, H., Lee, C., Park, H., Moon, S.: What is twitter, a social network or a news media?. In: *WWW*, pp. 591–600. ACM (2010)
10. Lau, J.H., Collier, N., Baldwin, T.: On-line trend analysis with topic models: \#twitter trends detection topic model online. In: *COLING*, pp. 1519–1534 (2012)
11. Mathioudakis, M., Koudas, N.: Twittermonitor: trend detection over the twitter stream. In: *SIGMOD*, pp. 1155–1158. ACM (2010)
12. McCallum, A.K.: Mallet: A machine learning for language toolkit (2002)
13. Naaman, M., Becker, H., Gravano, L.: Hip and trendy: characterizing emerging trends on twitter. *J. Am. Soc. Inf. Sci. Technol.* **62**(5), 902–918 (2011)

14. Nanavati, A.A., Gurumurthy, S., Das, G., Chakraborty, D., Dasgupta, K., Mukherjea, S., Joshi, A.: On the structural properties of massive telecom call graphs: findings and implications. In: Proceedings of the 15th ACM International Conference On Information And Knowledge Management, pp. 435–444. ACM (2006)
15. Narang, K., Nagar, S., Mehta, S., Subramaniam, L.V., Dey, K.: Discovery and analysis of evolving topical social discussions on unstructured microblogs. In: Advances in Information Retrieval, pp. 545–556. Springer (2013)
16. Stilo, G., Velardi, P.: Hashtag sense clustering based on temporal similarity. *Comput. Linguist.* (2017)
17. Yang, J., Leskovec, J.: Patterns of temporal variation in online media. In: WSDM, pp. 177–186. ACM (2011)